

## **METHOD AND APPARATUS FOR EFFICIENTLY REPRESENTING, STORING AND ACCESSING VIDEO INFORMATION**

The invention claims benefit of U.S. Provisional Application Number  
5 60/031,003, filed November 15, 1996.

The invention relates to video processing techniques and, more particularly, the invention relates to a method and apparatus for efficiently storing and accessing video information.

### 10 BACKGROUND OF THE DISCLOSURE

The capturing of analog video signals in the consumer, industrial and government/military environments is well known. For example, a moderately priced personal computer including a video capture board is typically capable of converting an analog video input signal into a digital video signal, and  
15 storing the digital video signal in a mass storage device (e.g., a hard disk drive). However, the usefulness of the stored digital video signal is limited due to the sequential nature of present video access techniques. These techniques treat the stored video information as merely a digital representation of a sequential analog information stream. That is, stored video is accessed in a  
20 linear manner using familiar VCR-like commands, such as the PLAY, STOP, FAST FORWARD, REWIND and the like. Moreover, a lack of annotation and manipulation tools due to, e.g., the enormous amount of data inherent in a video signal, precludes the use of rapid access and manipulation techniques common in database management applications.

25 Therefore, a need exists in the art for a method and apparatus for analyzing and annotating raw video information to produce a video information database having properties that facilitate a plurality of non-linear access techniques.

### 30 SUMMARY OF THE INVENTION

The invention is a method and apparatus for comprehensively representing video information in a manner facilitating indexing of the video information. Specifically, a method according to the invention comprises the steps of dividing a continuous video stream into a plurality of video scenes;

and at least one of the steps of dividing, using intra-scene motion analysis, at least one of the plurality of scenes into one or more layers; representing, as a mosaic, at least one of the plurality of scenes; computing, for at least one layer or scene, one or more content-related appearance attributes; and storing, in a database, the content-related appearance attributes or said mosaic representations.

### BRIEF DESCRIPTION OF THE DRAWINGS

The teachings of the present invention can be readily understood by considering the following detailed description in conjunction with the accompanying drawings, in which:

FIG. 1 depicts a high level block diagram of a video information processing system according to the invention;

FIG. 2 is a flow diagram of a segmentation routine suitable for use in the video information processing system of FIG. 1;

FIG. 3 is a flow diagram of an authoring routine suitable for use in the video information processing system of FIG. 1;

FIG. 4 depicts a "Video-Map" embodiment of the invention suitable for use as a stand-alone system, or as a client within the video information processing system of FIG. 1;

FIG. 5 shows a user holding the Video-Map embodiment of FIG. 4, and an exemplary screen display of an annotated image of the skyline of New York city;

FIG. 6 depicts exemplary implementation and use steps of the Video-Map embodiment of FIG. 4; and

FIG. 7 is a graphical representation of the relative memory requirements of two scene storage methods.

FIG. 8 is a flow diagram of a query execution routine according to the invention; and

FIGS. 9 and 10 are, respectively, a flow diagram 900 and a high-level function diagram 1000 of an attribute generation method according to the invention.

5       The invention will be described within the context of a video information  
processing system. It will be recognized by those skilled in the art that  
various other embodiments of the invention may be realized using the  
teachings of the following description. As examples of such embodiments, a  
video-on-demand embodiment and a "Video-Map" embodiment will also be  
10 described.

Segmenting comprises the process of dividing a continuous video stream into a plurality of segments, or scenes, where each scene comprises a plurality of frames, one of which is designated a “key frame.”

Mosaic construction comprises the process of computing, for a given scene or video segment, a variety of “mosaic” representations and associated frame coordinate transforms, such as background mosaics, synopsis mosaics, depth layers, parallax maps, frame-mosaic coordinate transforms, and frame-reference image coordinate transforms.. For example, in one mosaic representation a single mosaic is constructed to represent the background scenery in a scene, while individual frames in the scene include only foreground information that is related to the mosaic by an affine or a

5

15

ancillary data streams (time, sensor data, telemetry) or manual entry, ancillary data related to some or all of the scenes or video segments.

20

30

content may be used as guides in determining the appropriate representation level of the scene.

Fig. 1 is a high level block diagram of a video information processing system 100 according to the invention. The video information processing system 100 comprises three functional subsystems, an authoring sub-system, an access sub-system and a distribution sub-system. The three functional subsystems non-exclusively utilize various functional blocks within the video information processing system 100. Each of the three sub-systems will be described in more detail below, and with respect to the various drawings. Briefly, the authoring sub-system 120, 140 is used to generate and store a representation of pertinent aspects of raw video information and, specifically, to logically segment, analyze and efficiently represent raw video information to produce a video information database having properties that facilitate a plurality of access techniques. The access sub-system 130, 125, 150 is used to access the video information database according access techniques such as textual or visual indexing and attribute query techniques, dynamic browsing techniques and other iterative and relational information retrieval techniques. The distribution sub-system 130, 160, 170 is used to process accessed video information to produce video information streams having properties that facilitate controllably accurate or appropriate information stream retrieval and compositing by a client. Client-side compositing comprises the steps necessary to retrieve specific information in a form sufficient to achieve a client-side purpose.

Video information processing system 100 receives a video signal S1 from a video signal source (not shown). The video signal S1 is coupled to an authoring sub-system 120 and an image vault 150. The authoring subsystem 120 processes the video signal S1 to produce a video information database 125 having properties that facilitate a plurality of access techniques. For example, the video representative information resulting from the previously-mentioned comprehensive representation steps (i.e., segmenting, mosaic construction, motion analysis, appearance analysis and ancillary data capture) is stored in video information database 125. Video information database 125, in response to a control C1 requesting, e.g., video frames or scenes substantially matching some or all of the stored video representative

08970889 "1149"

5

15

20

25

30

The authoring and access subsystems will first be described in a general manner with respect to the video information processing system 100 of FIG. 1. The distribution subsystem will then be described within the context of several embodiments of the invention. In describing the several  
5   embodiments of the invention, several differences in the implementation of the authoring and access subsystems with respect to the embodiments will be noted.

The inventors have recognized that the problems of video sequence segmentation and video sequence searching may be addressed by the use of a  
10   short, yet highly representative description of the contents of the images. This description is in the form of a low-dimensional vector of real-valued quantities defined by the inventors as a multi-dimensional feature vector (MDFV). The MDFV "descriptor" comprises a vector descriptor of a predetermined dimensionality that is representative of one or more attributes  
15   associated with an image. An MDFV is generated by subjecting an image to a predetermined set of digital filters, where each filter is tuned to a specific range of spatial frequencies and orientations. The filters, when taken together, cover a wide range of spatial-frequencies and orientations. The respective output signals from the filters are converted into an energy  
20   representation by, e.g., integrating the squared modulus of the filtered image over the image region. The MDFV comprises these energy measures.

FIGS. 9 and 10 are, respectively, a flow diagram 900 and a high-level function diagram 1000 of an attribute generation method according to the invention. The method of FIG. 9 will be described with reference to FIG. 10.  
25   Specifically, the method 900 and function diagram 1000 are directed toward the processing of an input image  $I_0$  to produce attribute information (i.e., MDFV<sub>0</sub>) in the form of an attribute pyramid.

For the purposes of appearance-based indexing, two kinds of multi-dimensional features are computed: (1) Features that capture distributions  
30   without capturing any spatial constraints; and (2) Features that compute local appearance and are grouped together to capture the global spatial arrangement.

The first type of features that are computed do not preserve the spatial arrangement of the features within a layer or object. As described

46970889 1149

15

20

30

Multi-dimensional feature vectors are next computed at each salient location. That is, filter responses for filters at multiple scales and orientations



are computed. These may be defined using Gaussian and derivative filters or Gabor filters. A collection of these filters that systematically sample the space of orientations and scales (within reasonable limits, for instance scale changes between 1/8 and 8, but in principle may be arbitrary) is computed.

- 5 This collection as each of the salient points becomes the multi-dimensional feature representation for that point. For each layer and object, a collection of these features along with their spatial locations is stored in a database using a kd-tree (R-tree) like multi-dimensional data structure.

- The attribute generation method 900 of FIG. 9 is entered at step 905, 10 when an input frame is made available. At step 910 the input frame is retrieved, and at step 915 the input frame is subjected to a known pyramid processing step (e.g., decimation) to produce an image pyramid. In FIG. 10, the input frame is depicted as an input image  $I_0$ , and the pyramid processing step produces an image pyramid comprising three image pyramid subbands, 15  $I_1$ ,  $I_2$  and  $I_3$ .  $I_1$  is produced by, e.g., subsampling  $I_0$ .  $I_2$  is produced by, e.g., subsampling  $I_1$ .  $I_3$  is produced by, e.g., subsampling  $I_2$ . Since each subband of the image pyramid will be processed in the same manner, only the processing of subband  $I_1$  will be described in detail. Moreover, an image pyramid comprising any number of subbands may be used. A suitable pyramid 20 generation method is described in commonly assigned and copending U.S. Application 08/511,258, entitled METHOD AND APPARATUS FOR GENERATING IMAGE TEXTURES, filed August 4, 1995, and incorporated herein by reference in its entirety.

- After generating an image pyramid (step 915) the attribute generation 25 method 900 of FIG. 9 proceeds to step 920, where an attribute feature and an associated filtering scheme are selected, and to step 925, where  $N$  feature filter are used to filter each of the subbands of the image pyramid. In FIG. 10 the image subband  $I_1$  is coupled to a digital filter  $F_1$  comprising three subfilters  $f_1$ - $f_3$ . Each of the three subfilters is tuned to a specific, narrow range 30 of spatial frequencies and orientations. The type of filtering used, the number of filters used, and the range of each filter is adjusted to emphasis the type of attribute information produced. For example, the inventors have determined that color attributes are appropriately emphasized by using Gaussian filters, while texture attributes are appropriately emphasized by using oriented

08970889 "11497

filters (i.e., filters looking for contrast information in differing pixel orientations). It must be noted that more or less than three sub-filters may be used, and that the filters may be of different types.

After filtering each of the image pyramid subbands (step 925), the attribute generation method 900 of FIG. 9 proceeds to step 930, where the filter output signals are rectified to remove any negative components. In FIG. 10, the output signal from each of the three subfilters  $f_1$ - $f_3$  of digital filter  $F_1$  is coupled to a respective subrectifier within a rectifier  $R_1$ . The rectifier  $R_1$  removes negative terms by, e.g., squaring the respective output signals.

After rectifying each of the filter output signals (step 930), the attribute generation method 900 of FIG. 9 proceeds to step 935, where a feature map is generated for the attributes represented by each rectified filter output signal. In FIG. 10, feature map  $FM_1$  comprises three feature maps associated with, e.g., three spatial frequencies and orientations of subband image  $I_1$ . The three feature maps are then integrated to produce a single attribute representation  $FM_1$  of subband image  $I_1$ .

After generating the feature maps (step 935), the attribute generation method 900 of FIG. 9 proceeds to step 940, where the respective feature maps of each subband are integrated together in one or more integration operations to produce an attribute pyramid. In FIG. 10, the previously-described processing of subband image  $I_1$  is performed for subband images  $I_2$  and  $I_3$  in substantially the same manner.

After producing the attribute pyramid related to a particular attribute (step 940), the routine 900 of FIG. 9 proceeds to step 945, where the attribute pyramid is stored, and to step 945, where a query is made as to whether any additional features of the image pyramid are to be examined. If the query at step 945 is affirmatively answered, then the routine 900 proceeds to step 920, where the next feature and its associated filter are selected. Steps 925-950 are then repeated. If the query at step 945 is negatively answered, then the routine 900 proceeds to step 955, where a query is made as to whether the next frame should be processed. If the query at step 955 is affirmatively answered, then the routine 900 proceeds to step 910, where the next frame is input. Steps 915-955 are then repeated. If the

query at step 955 is negatively answered, then the routine 900 exits at step 960.

It is important to note that the attribute information generated using the above-described attribute generation method 900, 1000 occupies much less memory space than the video frame itself. Moreover, a plurality of such attributes stored in non-pyramid or pyramid form comprise an index to the underlying video information that may be efficiently accessed and searched, as will be described below.

The first functional subsystem of the video information processing system 100 of FIG. 1, the authoring sub-system 120, will now be described in detail. As previously noted, the authoring sub-system 120 is used to generate and store a representation of pertinent aspects of raw video information, such as information present in video signal S1. In the information processing system 100 of FIG. 1, the authoring subsystem 120 is implemented using three functional blocks, a video segmentor 122, an analysis engine 124 and a video information database 125. Specifically, the video segmentor 122 segments the video signal S1 into a plurality of logical segments, such as scenes, to produce a segmented video signal S2, including scene cut indicia. The analysis engine 124 analyzes one or more of a plurality of video information frames included within each segment (i.e., scene) in the segmented video signal S2 to produce an information stream S3. The information stream S3 couples, to an information database 125, information components generated by the analysis engine 124 that are used in the construction of the video information database 125. The video information database 125 may also include various annotations to the stored video information and ancillary information.

The segmentation, or “scene cut” function of the authoring subsystem 120 will now be described in detail. Video segmentation requires the detection of segment or scene boundaries using e.g., a “scene cut detector” that detects inter-frame discontinuities indicative of a change in scene, rather than a change in intra-frame information. This technique utilizes the fact that consecutive video frames are highly correlated and, in most cases, all frames in a particular scene have many attributes in common. A common example to an attribute used for scene cut detection is the background. Each scene

shot is assumed to have a single background and was shot at a single location, possibly from a small range of camera viewpoints.

FIG. 2 is a flow diagram of a segmentation routine suitable for use in the video information processing system of FIG. 1.

5       The segmentation routine 200 is entered at step 205, when the first frame of a new scene is received. The segmentation routine 200 then proceeds to step 210, where an index variable N is initialized to 1, and to step 220, where at least one of the above-described vector descriptors are calculated for the Nth frame. The segmentation routine 200  
10 then proceeds to step 230, where vector descriptors corresponding to those calculated at step 220 are calculated for the Nth+1 frame. Steps 220 and 230 may be implemented according to the principles of the attribute generation routine 900 discussed above.

After calculating the representative MDFV descriptors for the Nth  
15 (step 220) and Nth+1 (step 230) frames, the segmentation routine 200 then proceeds to step 235, where the difference (e.g., the Euclidian distance) between the Nth and Nth+1 MDFV descriptors is computed, to produce an interframe feature distance (IFFD). The segmentation routine 200 then proceeds to step 240, where the IFFD is compared to a threshold level. If the  
20 exceeds the threshold level (i.e., frame N is different than frame N+1 by the threshold amount), then the segmentation routine 200 proceeds to step 250, where the scene cut flag is set, and to step 255, where the segmentation routine 200 is exited. If the IFFD does not exceed the threshold level, then the index variable N is incremented by one (step 245), and steps 225-240 are  
25 repeated until a scene cut is detected.

The IFFD threshold level is either a predetermined level or, preferably, computed using the IFFD statistics of the available frames. Typically, this threshold is related to a "median" or other rank measures of the input set (i.e., the MDFV descriptors of the input frames). The segmentation routine 200 is  
30 depicted as operating in a single pass mode. However, the segmentation routine 200 can also be implemented in a two pass mode. In the single pass mode, the IFFD threshold level statistics are preferably determined on a "running" basis (e.g., a rolling average or other statistic based on the M most recent frames). In the two-pass mode, the IFFD threshold level statistics are

08970889 " 1149  
" 68802680



5

10

20

25

After the scene is optionally segmented into background and foreground portions, the routine 300 proceeds to step 315, where intra-scene attributes (i.e., intra-segment or frame-to-frame attributes) of each scene in

the segmented video information stream S2 are calculated. Intra-scene attributes, which will be discussed in more detail below, comprise intra-frame and inter-frame attributes of video frames within a particular video scene (i.e., attributes characteristic of one or more of the video information frames forming a scene). The previously described multi-dimensional feature vectors (MDFV<sub>s</sub>) may be used as intra-scene attributes. The analysis routine 300 then proceeds to step 320, where the calculated intra-scene attributes are stored in a video attribute database, such as video information database 125.

After calculating the intra-scene attributes of each scene, the analysis routine 300 proceeds to step 325, where inter-scene attributes (i.e., inter-segment or scene-to-scene attributes) of the segmented video information stream S2 are calculated. Inter-scene attributes, which will be discussed in more detail below, comprise attributes characteristic of one or more of the scenes forming a group of scenes (e.g., temporal order and the like). The calculation of step 325 utilizes information generated at step 315 and other information. The analysis routine 300 then proceeds to step 330, where the calculated inter-scene attributes are stored in a video attribute database, such as video information database 125.

After calculating the inter-scene attributes of the segmented video information stream S2, the analysis routine 300 then proceeds to optional step 335, where inter-scene representations, or "groupings" are calculated. The analysis routine 300 then proceeds to optional step 340, where the calculated representations are stored in a video attribute database, such as video information database 125. Inter-scene representations, which will be discussed in more detail below, comprise logical groupings of scenes to produce expanded visual representations of common subject matter (e.g., mosaics, 3D models and the like). The inter-scene grouping calculation and storage steps are optional because such representations or groupings are not used in every application.

The analysis routine 300 exits at step 345 when the input video signal S1 has been fully processed by the various functional blocks of the authoring sub-system. The result of the analysis routine 300 is a video attribute database, such as video information database 125, that includes a plethora of information related to the input video signal S1.

In the video information processing system 100 of FIG. 1, the input video signal S1, in a compressed or uncompressed form, is stored in image vault 150. Since one of the attributes of a scene is the presentation time of the scene (i.e., the time relative to the start of the video program that includes the scene), a scene identified using the video information database 125 may be retrieved from the image vault by retrieving the video information having the same presentation time.

The above-described analysis routine 300 refers to intra-scene attributes, inter-scene attributes, and inter-scene groupings. These concepts will now be described in detail.

Video information comprises a sequence or collection of video information frames, where each video frame is associated with a set of attributes. The set of attributes associated with a particular frame may be classified in a number of ways. For example, frame-specific attributes are those attributes of a video information frame that relate to the arrangement of video information within the particular frame. Examples of frame-specific attributes include distributions of luminance, chrominance, texture and shape; location coordinates of objects; textual and visual annotations and descriptions and the like. Segment-specific attributes are those attributes of a video information frame that relate to the arrangement of video information within a segment, or scene, comprising a plurality of video information frames. Examples of segment-specific attributes include the frame number of a particular video frame in a sequence of video frames, identification of a scene that the particular video frame is part of, geographic location and temporal information relating to the scene, static and dynamic geometric information relating to camera location(s) and usage(e.g., parallax information), identification of actors and objects within the scene, and the like. Other classifications may also be used, several of which will be discussed in other portions of this disclosure. Moreover, individual attributes may be utilized within a number of classifications.

In addition to intra-scene or intra-segment attributes, such as the frame-specific and segment-specific attributes derived directly from respective frame parameters and segment parameters, collections of frames or segments (sequential or otherwise) may be associated with "summaries,"

The above-described attribute classifications are used to generate a video information database 125 having properties that facilitate a plurality of access techniques. The video information database 125 will typically include intra-frame, inter-frame and inter-scene attribute data, any associated annotations, and address indicia associating the frame and scene attribute information with the actual video frames and scenes stored in the image vault 150. While the image vault 150 and the video information database 125 may be within the same mass storage device, this is not necessary. By accessing the attribute information using one or more of the various attribute classification sets, a user may access the video information frames and segments associated with attribute information. The user may also retrieve the stored attribute classification sets with or without the associated video information frames and segments, such as geometric information, dynamic information, ancillary information and the like.

It should be noted that it is not necessary to compute appearance  
 30 attributes for every frame in a particular scene, since such frames tend to be  
 highly correlated to begin with. Thus, the appearance attributes computed at  
 step 315 of the analysis routine 300 are computed only for “representative  
 frames,” e.g., mosaics or key frames within a scene. The selection of key  
 frames can be done automatically or manually for the specific application at



hand. Similarly, appearance attributes are computed for objects of interest, which may be defined either automatically using segmentation methods such as motion based segmentation, into coherently moving layers, or through color and texture analysis, or through manual outlining and specification of  
5 patches within a scene.

Appearance attributes of each representative frame and each object within a scene are computed independently and associated with the scene for subsequent indexing and retrieval of, e.g., the stored video. The appearance attributes consist of color and texture distributions, shape descriptions, and  
10 compact representations in terms of outputs of multiple scale, multiple orientation and multiple moment Gaussian and Gabor like filters. These attributes are organized in terms of data structures that will allow similarity queries to be answered very efficiently. For example, multi-dimensional R-tree data structures can be used for this purpose.

Each frame or scene in a video stream may be registered to a reference coordinate system. The reference coordinates are then stored along with the original video. This registration, or representation, of scenes allows, e.g., efficient storage of the video information comprising the scenes.  
15

After calculating the attribute information associated with the scenes comprising a particular program, the scenes may be grouped together and represented using one or more of a plurality of representation techniques. For example, video scenes may be represented using, e.g., two-dimensional mosaics, three-dimensional mosaics and networks of mosaics. A mosaic comprises an association, or joining, of a plurality of related video images to  
20 produce a combined video image having, e.g., additional field of view, panoramic effects and the like. In addition to providing new viewing experiences to a user, such representations of video information allow more efficient storage of the video information.

An example of a two-dimensional (2D) mosaic video representation is  
30 described in commonly assigned and copending U.S. Application No. 08/339,491 entitled SYSTEM FOR AUTOMATICALLY ALIGNING IMAGES TO FORM A MOSAIC IMAGE, filed November 14, 1994, and incorporated herein by reference in its entirety. In such a mosaic-based representation technique, a single mosaic is constructed to represent the

089089 " 1149  
4647T 68802680

background scenery in each scene. Each frame in the scene is related to the mosaic by an affine or a projective transformation. Thus, the 2D mosaic representation efficiently utilizes memory by storing the background information of a scene only once.

5           An example of a three-dimensional (3D) mosaic video representation is described in commonly assigned and copending U.S. Application No. 08/493,632, entitled METHOD AND SYSTEM FOR IMAGE COMBINATION USING A PARALLAX-BASED TECHNIQUE, filed June 22, 1995, and incorporated herein by reference in its entirety. A three-  
10 dimensional mosaic comprises a 2D image mosaic and a parallax mosaic. The parallax mosaic encodes the 3D structure of the scene. Each frame in the scene is related to the 3D mosaic by a 12 dimensional perspective transformation.

An example of a network of mosaics video representation is described in commonly assigned and copending U.S. Application No. 08/499,934, entitled METHOD AND SYSTEM FOR RENDERING AND COMBINING IMAGES, filed July 10, 1996, and incorporated herein by reference in its entirety. The network of mosaics comprises a network of 2D mosaics, where each mosaic corresponds to a single location. Each mosaic is constructed from the video captured by only rotating the camera about that single location. All mosaics are related to each other by coordinate transforms between them.

Video scenes may also be used to create three-dimensional structure models of various objects or portions of a scene. An interactive method to  
25 create a 3D structure model from video scenes is described in:

**"Reconstructing Polyhedral Models of Architectural Scenes from Photographs", C.J. Taylor, P.E. Debevec, and J. Malik, Proc. 4th European Conference on Computer Vision, Cambridge, UK, April 1996, pp. 659-668, incorporated herein by reference in its entirety.**

30 Video scenes may also be represented in terms of foreground and background. The above-incorporated U.S. Application No. 08/339,491 describes a technique for generating a model of the background portions of a scene. Foreground objects within the scene are obtained by aligning the background model with a video frame, and then subtracting the background

Video scenes may also be represented in terms of "layers." Layers are an extension to the basic mosaic concept for representing background motion. In the layered video representation, a separate mosaic "layer" is constructed for a foreground object. The foreground object is then tracked on a frame to frame basis by tracking the layer incorporating the object. Each shot is stored as a set of layered mosaics, a set of warping parameters for each layer for each frame, and a set of foreground residuals (if present). Representation of shots into layers may be achieved by techniques described in: "Layered Representation of Motion Video using Robust Maximum-Likelihood Estimation of Mixture Models and MDL Encoding", S. Ayer and H. Sawhney, Proc. IEEE Intl. Conference on Computer Vision, Cambridge, MA, June 1995, pp. 777-784, and in: "Accurate Computation of Optical Flow by using Layered Motion Representation", Proc. Intl. Conference on Pattern Recognition, Oct. 1994, pp. 743-746, both of which are incorporated by reference in their entireties.

20           The above-referenced layering techniques may be used in optional step  
310 of the analysis routine 300.

Scene representations, such as the mosaics, or other representations constructed for each frame, are grouped using their attributes to create a unified representation for all the frames. Since a movie or a sports event is typically imaged using just a few cameras and set locations, a large number of the frames will have a similar background. A possible criterion for grouping shots can therefore be a common background. In this case only one background mosaic needs to be stored for the entire group of frames. The grouping may be done manually, or automatically using techniques from the field of pattern recognition.

An automatic technique for grouping together scene shots based on color histograms was described in "Efficient Matching and Clustering of Video Shots", M. Yeung and B. Liu, IEEE Int. Conf. Image Processing, October 1995, Vol. A, pp. 338-341, incorporated herein by reference in its entirety.

To summarize, visual information is represented by a collection of scenes, or frame sequences. Each frame sequence typically includes a set of background and foreground models (e.g. mosaics), a viewing transformation which relates each frame to the appropriate model, and residual values for  
 5 each frame that correct for those effects that can not be represented by the model and viewing transformation. In addition to the visual information stored in, e.g., the image vault 150, appearance information related to the visual information is generated and stored in, e.g., video information database 125. Annotations, such as street names and various geographic, temporal and  
 10 relational data may also be stored in the database.

FIG. 7 is a graphical representation of the relative memory requirements of two scene storage methods. Specifically, structure and memory contents of a two-dimensional mosaic representation of a scene. A video program 710 comprises a plurality of scenes denoted as  $S_1$  through  $S_n$ .  
 15 A scene 720, illustratively scene  $S_{n-1}$ , comprises a plurality of video frames denoted as  $F_1$  through  $F_m$ , where  $F_1$  is the most recent frame. The video content of frames  $F_1$  and  $F_m$  are shown as respective pictures 730 and 740. Note that both pictures include a boat 732, 742 floating in a body of water 738, 748 beneath at least a portion of a cloud cover 736, 746. Picture 730  
 20 also contains a dock 739, while picture 740 contains the sun 744 but not the dock 739. Frames  $F_2$  through  $F_{m-1}$  are the intervening frames of scene 720, and show the scene changing from frame  $F_1$  to frame  $F_m$ .

Frame sequence 750 represents a two-dimensional mosaic of scene  $S_{n-1}$ . As previously discussed, the two-dimensional mosaic comprises a  
 25 background image related to all the frames in a particular scene, and a plurality of foreground images related to respective foreground portions of each frame of the scene. Thus, background frame 760 is shown as a panoramic picture comprising all the background information in scene  $S_{n-1}$ , i.e., a dock 769, a body of water 768, a cloud 766 and the sun 764. Frames  $F_1$   
 30 and  $F_m$  show only the respective foreground portions, comprising the boat (732, 742).

The pictures 730-740, and 760-780 are depicted in a graphical manner only for the purpose of illustrating the relative informational requirements to store each frame. It must be remembered that frames 770 and 780 require

"TTTTT" 6880/680

informational requirements for storing the standard frame sequence 720 of scene  $S_{n-1}$ , since the background portion of the scene, i.e., picture 760, is only stored once. Each of the frames within the two-dimensional mosaic representation of scene  $S_{n-1}$ , i.e., each of the frames within frame sequence 750, comprise only foreground and transform coordinate information.

15        Assuming that a video stream has been previously divided into  
subsequences, the access subsystem addresses, for example, the problem of  
finding a subsequence(s) to which a given frame belongs. This need arises  
during indexing and retrieval of stored video information for video editing and  
manipulation purposes. For example, given a representative frame from one  
20 subsequence, the user may be interested in determining other subsequences  
that contain images of the same scene.

The access sub-system is used to access the video information database using, textual query techniques, non-linear video browsing (i.e., "hyper-video") techniques, and linear browsing techniques. A textual query  
25 may comprise, e.g., a command to "find all video frames in a specific movie showing a specific actor," or "find all the touchdown scenes in all games that were played in a specific city during a specific period." A non-linear video browsing technique may comprise, e.g., iteratively grouping attribute-related video frames and video segments, where each successive frame or segment  
30 selection retrieves more relevant, or desirable, video information frame or segments for display. A linear video browsing technique may comprise, e.g., pointing to a specific displayed object, such as a baseball player, using a pointing device; and retrieving other scenes including the identified object (player), or displaying a list of all games played by this player. An object

Referring to FIG. 1, the access engine 130, in response to a textual, non-linear or linear access request from a user (e.g., from a client 170 via the network 160), accesses the video information database and identifies video frames and/or scenes together with geometric, dynamic and other scene structure information that satisfy the user request. As previously noted, the video information database 125 will typically include intra-frame, inter-frame and inter-scene attribute data, any associated annotations, and address indicia associating the frame and scene attribute information with the actual video frames and scenes stored in the image vault 150. A user may interactively access the attribute data alone, or the attribute data in conjunction with the actual video frames and/or scenes. If the user wishes to view the actual video frames and/or scenes, then the access engine causes the image vault 150 to generate a video output signal S5. The video output signal S5 is then coupled to the user as signal S6.

The access engine 130 is capable of retrieving specific video information on a frame-by-frame basis by performing a search on a representative feature(s) of the desired video frames. As previously discussed, individual video frames are represented by a plurality of attributes which have been stored in a video information database 125. The access engine 130 utilizes the video information database 125 to retrieve, e.g., address indicia of frames or scenes corresponding to one or more desired attributes.

FIG. 8 is a flow diagram of a query execution routine according to the invention. A methodology for searching for individual video frames in the available frame subsequences (i.e., scenes) relies on the previously-described multi-dimensional feature vector descriptor representation of individual frames, and assumes that the input sequence has been previously segmented into subsequences and processed by the authoring subsystem 120.

The routine 800 is entered when a user specifies a query type (step 805) and a query specification (step 810). A query type comprises, e.g., color, texture, keywords and the like. A query specification is a more particular

5

10

20

25

30



5

10

25

30



particular location in an image (or input of map, GPS or other reference coordinates), video clips associated with that location may then be accessed.

In the case of, e.g., a mosaic representation video information having desired attributes, the access subsystem uses the transformation between  
5 the video frames and the image representation to retrieve other frames or scenes where the particular location or attribute is visible. This technique is described in commonly owned and copending U.S. Application No. 08/663,582 entitled A SYSTEM FOR INDEXING AND EDITING VIDEO SEQUENCES USING A GLOBAL REFERENCE filed June 14, 1996 and incorporated  
10 herein by reference in its entirety.

The presentation of video information, still image information and other information retrieved from the video information database 125 or the image vault 150 may be adapted to suit particular applications of the invention. For example, the presented information may be annotated or unannotated.  
15 Moreover, the presentation may be adapted to facilitate further querying. The following is a partial list of the video information presentation possibilities.

The video information may be presented as a single video frame, or a collection of isolated frames, in response to a user's query. Such frames are part of the original images and video sequences used to create the video  
20 information database. Similarly, the video information may be presented as a single scene, or a collection of scenes, from the original video. The video information may be presented in one of the previously described mosaic formats. Such a mosaic is usually pre-computed prior to a query, and is displayed, completely or in part, as an answer to the query.

25 The video information may be presented as one or more newly generated images. For example, when queried using positional information, the system can generate a new view of a scene or object as seen from that particular viewing position. Methods to use video representations to create a desired view are described in U.S. Application No. 08/493,632 and in U.S.  
30 Application No. 08/499,934. Other methods for new view generation, like those using a 3D CAD model, can be used as well. An example is described in "Reconstructing Polyhedral Models of Architectural Scenes from Photographs", C.J. Taylor, P.E. Debevec, and J. Malik, Proc. 4th European

0897088 "1149  
Z64TT" 68802680

Conference on Computer Vision, Cambridge, UK, April 1996, pp. 659-668, incorporated herein by reference in its entirety.

The video information may be presented in a manner that highlights dynamic content (e.g., foreground or moving objects). For example, in order to  
5 more clearly view moving objects and other dynamic content as well as the static background, the dynamic content can be overlaid on a static summary mosaic of the background to show a compete summary of the video in an expanded view format.

FIG. 4 depicts a "Video-Map" embodiment 470 of the invention suitable  
10 for use as a stand-alone system, or as a client 170-2 within the video information processing system 100 of FIG. 1. The Video-Map 470 comprises a display 472, a network interface 473, a controller 474 and an input device 475 that operate in substantially the same manner as previously described with respect to the client 170 of FIG. 1. The Video-Map 470 also includes one  
15 or more ancillary information sources 476 suitable for providing positioning information, illustratively a Global Positioning System (GPS) receiver 476-1 and a digital camera 476-2. The ancillary information source(s) 476 provide information that is used by the controller 474 to generate video information database queries.

20 The Video-Map 470 optionally includes a video storage unit 477, such as a CD-ROM drive, that is coupled to the controller 474 via a video storage unit interface 478. The video storage unit 477 is used to store an annotated video information database, such as the annotated video information database 125 similar to that of the information processing system 100 of  
25 FIG. 1. The video storage interface 478, in conjunction with the controller 474, performs substantially the same function as the access engine 130 of the information processing system 100 of FIG. 1.

The Video-Map 470, in the client mode of operation, communicates with the access engine 130 of the information processing system 100 via  
30 network interface 173, which is coupled to a network 160, illustratively a cellular or satellite telecommunications network 160.

The purpose of the Video-Map embodiment is to capture, annotate and represent visual and other information about a geographic environment in a structured form, and to be able to access and present both the visual and

08970889 111497  
264111 58802680

other information at a later time in a form that situates the browser in the geometric and visual context of the current environment.

FIG. 5 shows a user 505 holding the Video-Map embodiment 470 of FIG. 4, and an exemplary screen display 510 of an annotated image of the skyline of New York city. It should be noted that the displayed image is similar to what the user sees with his eyes. However, the displayed image is annotated such that many of the buildings are identified by corresponding text 521, 522, 523. The information necessary to produce the displayed image is stored in an annotated video information database either locally (i.e., in the video storage unit 472) or remotely (i.e., in the video information database 125 of FIG. 1).

The representation of the city of New York stored in the local or remote video information database includes the geometric, visual and ancillary information about landmarks and locales of interest. This annotated representation is created from video images captured through a variety of sources, and from mapping and ancillary information obtained from other sources. This annotated database is typically stored in a compressed format on one or more storage platforms. To conserve memory and processing resources, the displayed image may be a still image.

The stored database is accessed by providing ancillary information that approximately locates the user within the coordinate system space of the video information representation stored in the video information database. Such ancillary information may include positional data, e.g., data retrieved from the GPS receiver 476-1. The positional information forms the basis of a query into the video information database. That is, the controller 474 constructs a query of the form "show all portions of the New York city skyline visible from this location." In the client mode of operation, the query is transmitted to the access engine 130 via the network in the previously described manner. The access engine retrieves the appropriate view of New York City from the video information database 125, and coupled the retrieved view to the Video-Map 470 via the network 160. In the stand-alone mode of operation, the controller 474, in conjunction with the video storage interface 478, identifies and retrieves the appropriate view from the video storage unit

477. The appropriate view in either mode of operation may be coupled to the display 472 for viewing by the user.

The stored database is optionally accessed by providing ancillary information that includes single or multiple views in a visual form for the  
5 locale of interest, e.g., image data retrieved from the camera 476-2. The retrieved image data is subjected to an attribute identification process, and the resulting attribute information forms the basis of a query into the video information database.

In either the positional data case, or the visual attribute case, the  
10 access information is used to index into the video map database, and the retrieved information is presented to the viewer in a useful form. For example, the visual information may be presented in the form of an image/mosaic or video as would be seen from the viewpoint of the client. The presented information may optionally be annotated with textual, graphical or audible  
15 information, and other multi-modal annotations that are associated with the accessed locale. The annotations may be used to explain to the user the identity, function and other pre-stored relevant information of the objects in the presented field of view. Furthermore, the user may select, using the input device 475, different parts of the image to interactively access more  
20 information about a selected building or site of interest. The user can further query the system using any additional indices, such as hotel, restaurant, tourist interest and the like. Moreover, the Video-Map may be used as a navigation tool.

FIG. 6 depicts exemplary implementation and use steps of the Video-  
25 Map embodiment of FIG. 4. There are three main components of the video map embodiment of the invention: First, creating an annotated video map database (steps 610, 612, 613 and 614); Second, accessing the video map database (620, 622 and 624); and Third, presenting and viewing the visual and ancillary annotation information (630). It will be understood by those  
30 skilled in the art that the particular methods taught by this embodiment of the invention are not the only methods suitable for implementing the invention. Other methods useful to the practice of the invention are also contemplated to be within the scope of the invention. For example, in aerial

2025 RELEASE UNDER E.O. 14176

The first component of the Video-Map embodiment, creating an annotated video map database (i.e., authoring) will now be described. Starting with a collection of videos of a set of locales (e.g., New York), a video information database is generally constructed as previously described. The key to implementing the video map application is proper representation of the video information. Specifically, a collection of layered 2D and 3D mosaic images and parallax maps compactly represent the geometric and visual information of the locale (step 612). This representation of the actual video information is stored in the image vault 150 and video information database 125 or storage unit 477, along with the coordinate transforms that relate other such representations associated with a locale. The fundamental methodology for developing this representation was described above and in U.S. Application No. 08/493,632. This representation allows generation of either the original collection of videos that were used to create the representation, or new views of the same locales that were not present in any particular frame of the original video.

In addition to the representation of the geometric and visual information (step 612), two other classes of information are associated with the map database. One class represents the visual information not in terms of pixels and their color/intensity values (as is done in the above representation) but as higher order features that are computed from the pixel information. These features represent the distributions and spatial relationships of color, texture and shape like significant features of a locale that can describe the visual appearance of significant structures in a compact form (step 613). In general, these features are multidimensional vectors, matrices and tensors that encode significant visual appearances compactly. These features and their combinations will be used to index and match a specified query in the form of the appearance of an object/view of a locale at the time of map database access.

The third class of information associated with the map database consists of geographical map coordinates, GPS coordinates, textual descriptions of objects and views of a locale, audio/close-caption descriptions

5

information to access the relevant database locale.

10

20

30

5

10

15

25

30

15

25

30



video editing, managing and archiving large collections of videos for instance in government, military aerial video collections, and authoring multimedia content where videos are an important source of data. Therefore, the data representations, authoring tools and algorithms and user interaction and  
5 visualization tools may all be together or independently suited for a wide variety of video applications.

The information processing system 100 of FIG. 1 may be utilized as a video-on-demand (VOD) server. A client 170 in a VOD system will typically include a consumer television (i.e., display device 172), a remote control (i.e.,  
10 input device 175) and a set top terminal (i.e., controller 174 in combination with network interface 173). The VOD client-server application is directed to providing rapid program selection and program visualization to a client (i.e., subscriber).

Programs are stored in the image vault 150, and accessed by the  
15 access engine 130 in conjunction with the video information database 125. The database formation and access techniques are substantially the same as those techniques previously described. Additional access and distribution concerns involve billing and content restriction management.

The present invention can be embodied in the form of computer-  
20 implemented processes and apparatuses for practicing those processes. The present invention also can be embodied in the form of computer program code embodied in tangible media, such as floppy diskettes, CD-ROMs, hard drives, or any other computer readable storage medium, wherein , when the computer program code is loaded into and executed by a computer, the  
25 computer becomes an apparatus for practicing the invention. The present invention can also be embodied in the form of computer program code, for example whether stored in a storage medium, loaded into and/or executed by a computer, or transmitted over some transmission medium, such as over electrical wiring or cabling, through fiber optics, or via electromagnetic  
30 radiation, wherein, when the computer program code is loaded into and executed by a computer, the computer becomes an apparatus for practicing the invention. When implemented on a general-purpose microprocessor, the computer program code segments configure the microprocessor to create specific logic circuits.

08970889 "1149"

Although various embodiments which incorporate the teachings of the present invention have been shown and described in detail herein, those skilled in the art can readily devise many other varied embodiments that still incorporate these teachings, such as computer-implemented processes and  
5 apparatuses for practicing those processes.

03970889 11497  
" 68802680